9

10

11

1

2

3

2.

CLAIMS

What is claimed is:

1	1.	A method for distributing data items from a particular set of data into a plurality of
2		buckets based on distribution keys associated with said data items, the method
3		comprising the steps of:
4,		randomly selecting data items from said particular set of data to produce a

randomly selecting data items from said particular set of data to produce a sampled set of data items;

determining a plurality of ranges based on the distribution keys associated with the sampled set of data items;

assigning said plurality of ranges to said plurality of buckets; and distributing each data item in said particular set of data to the bucket that has been assigned the range into which falls the distribution key of the data item.

- The method of Claim 1 wherein the step of randomly selecting data items from said particular set of data includes randomly selecting data items from each subset of a plurality of subsets of said particular set of data.
- The method of Claim 2 wherein the step of randomly selecting data items from each subset of a plurality of subsets of said particular set of data includes randomly selecting data items from each partition of a partitioned table.

1
` ~ _
ands.
Ţ
Ħ
I
IJ
::
, die
=
بالد أ
, F

3

4

5

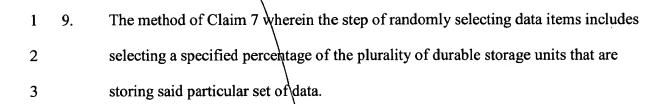
6

1	4.	The method of Claim 2 wherein the step of randomly selecting data items from each
2		subset of a plurality of subsets of said particular set of data includes randomly
3		selecting data items from subsets of data, stored in buffers in volatile memory, that
4		represent results of one or more previously performed operations.

- The method of Claim 1 further comprising the steps of:
- 2 assigning the plurality of buckets to a plurality of processes; and

1

- 3 causing each process of said plurality of processes to perform, in parallel with the other processes of said plurality of processes, an operation on the data items 4
- contained in any buckets assigned to the process. 5
- The method of Claim 2 further comprising the step of selecting a distinct random seed 1 6. for each subset of the plurality of subsets of said particular set of data. 2
- 7. The method of Claim 1 wherein: 1
 - the particular set of data is durably stored on a plurality of durable storage units; and the step of randomly selecting data items from said particular set of data to produce a sampled set of data items includes randomly selecting durable storage units from said plurality of durable storage units and using the data items stored on said randomly selected durable storage units as the sampled set of data items.
- The method of Claim 1 wherein the step of randomly selecting data items includes 1 8. selecting a specified percentage of data items in said particular set of data. 2



- 1 10. The method of Claim 8 further comprising the step of receiving, from a user, data that
 2 specifies said percentage.
- 1 11. The method of Claim 9 further comprising the step of receiving, from a user, data that
 2 specifies said percentage.
- The method of Claim 5 wherein said operation is specified in a database command, the method further comprising receiving with said database command data that indicates how much of said particular set of data to randomly select to produce said sampled set of data items.
- The method of Claim 1 wherein the step of determining a plurality of ranges based on the distribution keys associated with the sampled set of data items includes determining ranges that contain an approximately equal amount of distribution keys associated with said sampled set of data items.
- 1 14. A computer-readable medium carrying instructions for distributing data items from a
 2 particular set of data into a plurality of buckets based on distribution keys associated
 3 with said data items, the instructions comprising instructions for performing the steps
 4 of:

5		randomly selecting data items from said particular set of data to produce a
6		sampled set of data items;
7		determining a plurality of ranges based on the distribution keys associated
8		with the sampled set of data items;
9		assigning said plurality of ranges to said plurality of buckets; and
10		distributing each data item in said particular set of data to the bucket that has
11		been assigned the range into which falls the distribution key of the data
12		item.
1	15.	The computer-readable medium of Claim 14 wherein the step of randomly selecting
2		data items from said particular set of data includes randomly selecting data items from
3		each subset of a plurality of subsets of said particular set of data.
1	16.	The computer-readable medium of Claim 15 wherein the step of randomly selecting
2		data items from each subset of a plurality of subsets of said particular set of data
3		includes randomly selecting data items from each partition of a partitioned table.
1	17.	The computer-readable medium of Claim 15 wherein the step of randomly selecting
2		data items from each subset of a plurality of subsets of said particular set of data
3		includes randomly selecting data items from subsets of data, stored in buffers in
4.		volatile memory, that represent results of one or more previously performed
5		operations.
1	18.	The computer-readable medium of Claim 14 further comprising instructions for
2		performing the steps of:

2

3

4

5

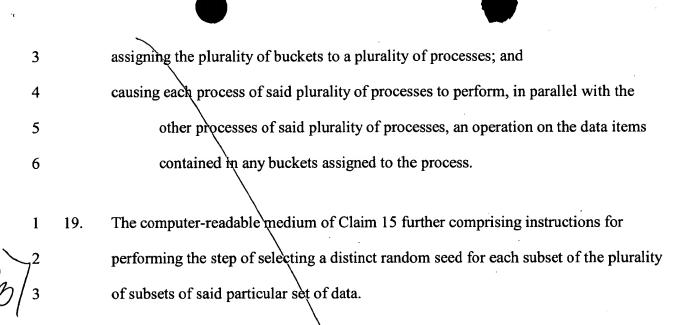
6

1

2

3

21.



20. The computer-readable medium of Claim 14 wherein:

the particular set of data is durably stored on a plurality of durable storage units; and the step of randomly selecting data items from said particular set of data to produce a sampled set of data items includes randomly selecting durable storage units from said plurality of durable storage units and using the data items stored on said randomly selected durable storage units as the sampled set of data items.

- The computer-readable medium of Claim 14 wherein the step of randomly selecting data items includes selecting a specified percentage of data items in said particular set of data.
- The computer-readable medium of Claim 20 wherein the step of randomly selecting data items includes selecting a specified percentage of the plurality of durable storage units that are storing said particular set of data.

2

3

4

1

2

3

4

26.

- 1 23. The computer-readable medium of Claim 21 further comprising instructions for 2 performing the step of receiving, from a user, data that specifies said percentage.
- The computer-readable medium of Claim 22 further comprising instructions for performing the step of receiving, from a user, data that specifies said percentage.
 - 25. The computer-readable medium of Claim 18 wherein said operation is specified in a database command, the computer-readable medium further comprising instructions for receiving with said database command data that indicates how much of said particular set of data to randomly select to produce said sampled set of data items.
 - The computer-readable medium of Claim 14 wherein the step of determining a plurality of ranges based on the distribution keys associated with the sampled set of data items includes determining ranges that contain an approximately equal amount of distribution keys associated with said sampled set of data items.